



Tema 2.2

TEORÍA DE LA ESTIMACIÓN

Febrero-Mayo 2006

1



INDICE

- 2.2.1 INTRODUCCIÓN
- 2.2.2 ESTIMACIÓN DE MÁXIMA VEROSIMILITUD (ML)
- 2.2.3 ESTIMACIÓN BAYESIANA
- 2.2.4 PROBLEMAS DE LA DIMENSIONALIDAD
 - ANÁLISIS DE COMPONENTES PRINCIPALES (PCA)
 - ANÁLISIS POR MÚLTIPLES DISCRIMINANTES (MDA)
- 2.2.5 CONCLUSIONES

2



2.2.1 INTRODUCCIÓN



La clasificación bayesiana precisa del conocimiento de $f_{\mathbf{x}}(\mathbf{x} | \omega_i)$ y de $\Pr(\omega_i)$. Para el cálculo de estas magnitudes se requiere:

- Disponer de una serie de datos previamente clasificados de forma fiable.
- Disponer de un estimador de esas probabilidades.

La estimación de $f_{\mathbf{x}}(\mathbf{x} | \omega_i)$ requiere muchos datos a menos que podamos definir una función que dependa de unos pocos parámetros θ_i .

Caso gaussiano: θ_i contiene la media y la matriz de covarianza

$$f_{\mathbf{x}}(\mathbf{x} | \omega_i, \theta_i) = \frac{1}{(2\pi)^{d/2} |\mathbf{C}_i|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \mathbf{C}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right\}$$

3



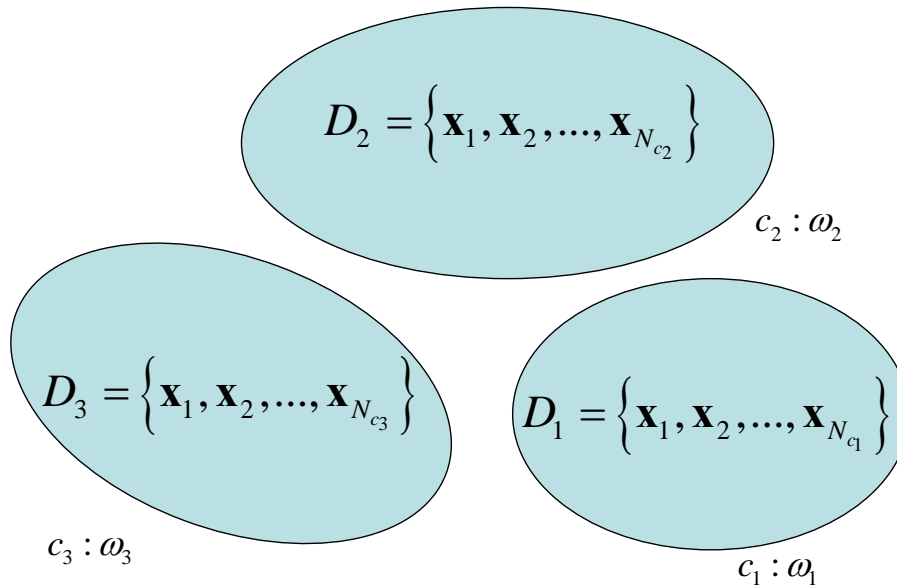
Existen dos alternativas:

1. **Estimación de máxima verosimilitud (ML):** Los parámetros a estimar se consideran deterministas (aunque desconocidos).
2. **Estimación bayesiana:** Los parámetros son variables aleatorias con distribución a priori conocida. La definición del estimador bayesiano permite mejorar fácilmente la estimación de $f_{\mathbf{x}}(\mathbf{x} | \omega_i)$ cuando se dispone de nuevos datos.

4



Supondremos que disponemos de una base de datos clasificada (un conjunto de vectores de características clasificados por categorías), a partir de las cuales hemos de determinar $f_{\mathbf{x}}(\mathbf{x}|\omega_i)$:



5



2.2.2 ESTIMACIÓN DE MÁXIMA VEROSIMILITUD (ML)

Si los datos observados en cada clase D son independientes:

$$f(D_i | \theta_i) = \prod_{k=1}^{N_{c_i}} f_{\mathbf{x}}(\mathbf{x}_k | \theta_i)$$

es la función de verosimilitud. El estimador maximiza esta función (o su logaritmo):

$$\hat{\theta}_{i,ML} = \arg \max_{\theta_i} f(D_i | \theta_i) = \arg \max_{\theta_i} \ln f(D_i | \theta_i)$$

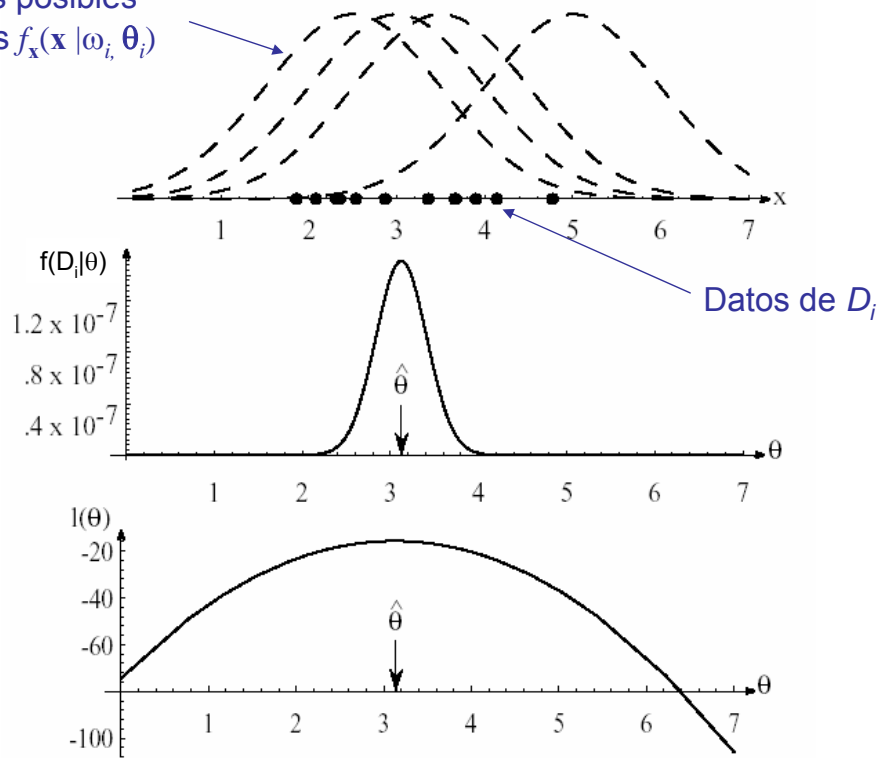
Un conjunto de condiciones necesarias (no suficientes) para obtener el estimador vienen dadas por:

$$\nabla_{\theta_i} \ln f(D_i | \theta_i) = \mathbf{0}$$

6



Algunas posibles funciones $f_x(x | \omega_i, \theta_i)$



7



Propiedades del estimador ML:

1. Es asintóticamente insesgado (en muchos casos es insesgado aunque N sea pequeño).
2. Es asintóticamente eficiente (cuando N es grande, su varianza es la de Crámer-Rao).

Sin embargo

1. No tiene porqué ser el que proporcione menor error de clasificación cuando utilizemos

$$f_x(x | \omega_i, \hat{\theta}_{i,ML})$$

2. Si la pdf asumida es muy distinta de la real las estimaciones pueden ser de poca calidad.

8



Ejemplo 1:

Estimador ML de la media μ_i si la matriz de covarianza C_i es conocida, en el caso gaussiano multivariable. Demostrad que:

$$\hat{\mu}_{i,ML} = \frac{1}{N_{c_i}} \sum_{k=1}^{N_{c_i}} \mathbf{x}_k$$

Ejemplo 2:

Estimador ML de la media μ_i y la matriz de covarianza C_i en el caso gaussiano multivariable. Demostrad que:

$$\hat{\mu}_{i,ML} = \frac{1}{N_{c_i}} \sum_{k=1}^{N_{c_i}} \mathbf{x}_k \quad \hat{C}_{i,ML} = \frac{1}{N_{c_i}} \sum_{k=1}^{N_{c_i}} (\mathbf{x}_k - \hat{\mu}_{i,ML})(\mathbf{x}_k - \hat{\mu}_{i,ML})^T$$

9



Ejemplo 3:

Estimador ML de la probabilidad p_k de aparición de '1' para cada una de las componentes del vector de datos binarios

$\mathbf{x} \in \{0,1\}^d$:

$$f_{\mathbf{x}}(D | \omega, \mathbf{p}) = \prod_{j=1}^{N_i} \prod_{k=1}^d p_k^{x_{k,j}} (1 - p_k)^{1-x_{k,j}}$$

$$\mathbf{p} = [p_1, \dots, p_d]$$

10



2.2.3 ESTIMACIÓN BAYESIANA



Si se dispone de algún conocimiento a priori de los rangos de valores más probables de θ_i podemos aprovecharlo para:

1. Mejorar la estimación ML de θ_i (usando MAP)

$$\hat{\theta}_{i,MAP} = \arg \max_{\theta_i} f(D_i | \theta_i) f(\theta_i) = \arg \max_{\theta_i} [\ln f(D_i | \theta_i) + \ln f(\theta_i)]$$

2. Estimar directamente las probabilidades a posteriori $\Pr(\omega_i | \mathbf{x})$

Calculando $f_{\mathbf{x}}(\mathbf{x} | \omega_i)$ y $\Pr(\omega_i)$. Es el procedimiento más aconsejable en una aplicación de clasificación.

11



Suposiciones

Queremos determinar la probabilidad a posteriori a partir de las observaciones en D_i , y supondremos que:

- La forma de $f_{\mathbf{x}}(\mathbf{x} | \theta_i)$ es conocida pero no el parámetro θ_i
- Nuestro conocimiento a priori de θ_i está en $f(\theta_i)$
- El resto de nuestro conocimiento sobre θ_i viene dado por los datos en D_i

12



Procedimiento:

1. Promediar la forma conocida para la función de verosimilitud respecto a la probabilidad a posteriori del parámetro:

$$f_{\mathbf{x}}(\mathbf{x} | \omega_i) \cong f_{\mathbf{x}}(\mathbf{x} | D_i) = \int f(\mathbf{x} | \boldsymbol{\theta}_i) f(\boldsymbol{\theta}_i | D_i) d\boldsymbol{\theta}_i$$

2. Calculamos la probabilidad a posteriori del parámetro como

$$f(\boldsymbol{\theta}_i | D_i) = \frac{f(D_i | \boldsymbol{\theta}_i) f(\boldsymbol{\theta}_i)}{\int f(D_i | \boldsymbol{\theta}_i) f(\boldsymbol{\theta}_i) d\boldsymbol{\theta}_i} \propto f(D_i | \boldsymbol{\theta}_i) f(\boldsymbol{\theta}_i)$$

3. Suponiendo independencia de los datos en D_i

$$f(D_i | \boldsymbol{\theta}_i) = \prod_{k=1}^{N_i} f(\mathbf{x}_k | \boldsymbol{\theta}_i)$$

13



Ejemplo 4:

Estimador bayesiano de $f_{\mathbf{x}}(\mathbf{x}|D_i)$ si

$$f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\mu}) \sim N(\boldsymbol{\mu}, \mathbf{C}) \quad f(\boldsymbol{\mu}) \sim N(\boldsymbol{\mu}_0, \mathbf{C}_0)$$

donde se suponen conocidas $\boldsymbol{\mu}_0$, \mathbf{C}_0 y \mathbf{C} , y se dispone de los datos observados $D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in \omega$

A partir de **2** y **3** podemos escribir:

$$\begin{aligned} f(\boldsymbol{\mu} | D) &= \alpha \prod_{k=1}^N f_{\mathbf{x}}(\mathbf{x}_k | \boldsymbol{\mu}) f(\boldsymbol{\mu}) = \\ &= \alpha' \exp \left[-\frac{1}{2} \left(\boldsymbol{\mu}^T (\mathbf{N}\mathbf{C}^{-1} + \mathbf{C}_0^{-1}) \boldsymbol{\mu} - 2\boldsymbol{\mu}^T \left(\mathbf{C}^{-1} \sum_{k=1}^N \mathbf{x}_k + \mathbf{C}_0^{-1} \boldsymbol{\mu}_0 \right) \right) \right] \end{aligned}$$

14



La ecuación puede escribirse también como:

$$f(\boldsymbol{\mu}|D) = \alpha'' \exp\left[-\frac{1}{2}(\boldsymbol{\mu} - \boldsymbol{\mu}_N)^T \mathbf{C}_N^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_N)\right]$$

donde:

$$\begin{aligned} \mathbf{C}_N^{-1} &= N\mathbf{C}^{-1} + \mathbf{C}_0^{-1} \\ \mathbf{C}_N^{-1}\boldsymbol{\mu}_N &= N\mathbf{C}^{-1}\mathbf{m}_N + \mathbf{C}_0^{-1}\boldsymbol{\mu}_0 \\ \mathbf{m}_N &= \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k \end{aligned}$$

Aplicando la igualdad matricial:

$$(\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} = \mathbf{A}(\mathbf{A} + \mathbf{B})^{-1} \mathbf{B} = \mathbf{B}(\mathbf{A} + \mathbf{B})^{-1} \mathbf{A}$$

15



pueden reformularse la media y la matriz de covarianza:

$$\begin{aligned} \boldsymbol{\mu}_N &= \mathbf{C}_0 \left(\mathbf{C}_0 + \frac{1}{N} \mathbf{C} \right)^{-1} \mathbf{m}_N + \frac{1}{N} \mathbf{C} \left(\mathbf{C}_0 + \frac{1}{N} \mathbf{C} \right)^{-1} \boldsymbol{\mu}_0 \\ \mathbf{C}_N &= \mathbf{C}_0 \left(\mathbf{C}_0 + \frac{1}{N} \mathbf{C} \right)^{-1} \frac{1}{N} \mathbf{C} \end{aligned}$$

Nótese que la media es una combinación lineal del conocimiento a priori de la media $\boldsymbol{\mu}_0$ y la información aportada por los datos \mathbf{m}_N . Integrando la ecuación 1:

$$f_{\mathbf{x}}(\mathbf{x} | \omega) \cong f_{\mathbf{x}}(\mathbf{x} | D) = \int f(\mathbf{x} | \boldsymbol{\mu}) f(\boldsymbol{\mu} | D) d\boldsymbol{\mu} \sim N(\boldsymbol{\mu}_N, \mathbf{C} + \mathbf{C}_N)$$

Cuando $N \rightarrow \infty$ la estimación tiende a ser ML

$$\boldsymbol{\mu}_N = \mathbf{m}_N \qquad \mathbf{C}_N = \frac{1}{N} \mathbf{C}$$

16



Comparación:

La función $f(D_i | \theta_i)$ tendrá un pico tanto más abrupto alrededor de $\theta_i = \hat{\theta}_i$ cuanto mayor sea N_i .

Si $f(\theta_i)$ no es cero y no varía mucho cerca de $\theta_i = \hat{\theta}_i$ entonces

$$f(\theta_i | D_i) = \frac{f(D_i | \theta_i)f(\theta_i)}{f(D_i)}$$

también tiene un pico en $\theta_i = \hat{\theta}_i$ y los estimadores obtenidos por Bayes y mediante ML coinciden. En la práctica, si el número de muestras de entrenamiento es pequeño, es mejor la estimación Bayesiana.



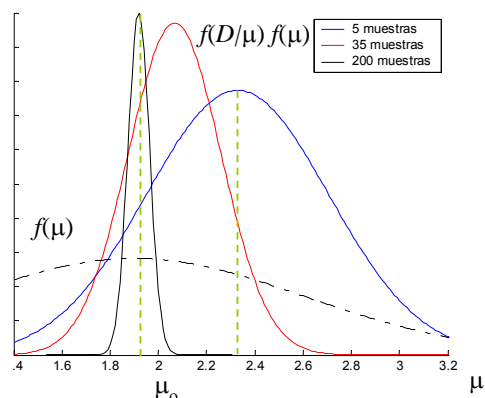
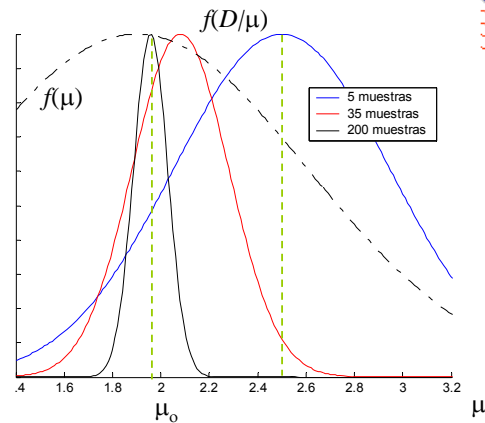
Ejemplo 5:

Estimación **ML** de la media ($\mu_0=2$) sobre un número variable de muestras Gaussianas.

La fdp a priori de μ es Gaussiana.

Estimación **Bayesiana** de la media ($\mu_0=2$) sobre un número variable de muestras Gaussianas.

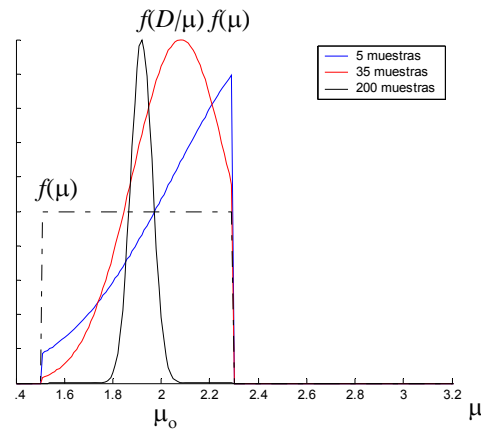
La fdp a priori de μ es Gaussiana.





Estimación **Bayesiana** de la media ($\mu_0=2$) sobre un número variable de muestras Gaussianas.

La fdp a priori de μ es uniforme.



19



2.2.4 PROBLEMAS DE DIMENSIONALIDAD

Al aumentar el número de características independientes en un problema de clasificación es posible hacer tender el error a cero.

Ejemplo 1:

Problema de clasificación con dos clases gaussianas equiprobables

$$f_{\mathbf{x}}(\mathbf{x} | \omega_i) \sim N(\boldsymbol{\mu}_i; \boldsymbol{\Sigma}) \quad i = 1, 2$$

$$\Pr(\omega_1) = \Pr(\omega_2)$$

$$\mathbf{x} \in \mathbb{R}^d$$

20



La probabilidad de error queda definida como:

$$\Pr(e) = \frac{1}{\sqrt{2\pi}} \int_{r/2}^{\infty} \exp(-u^2 / 2) du$$

$$r^2 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

Si las características son independientes entre sí:

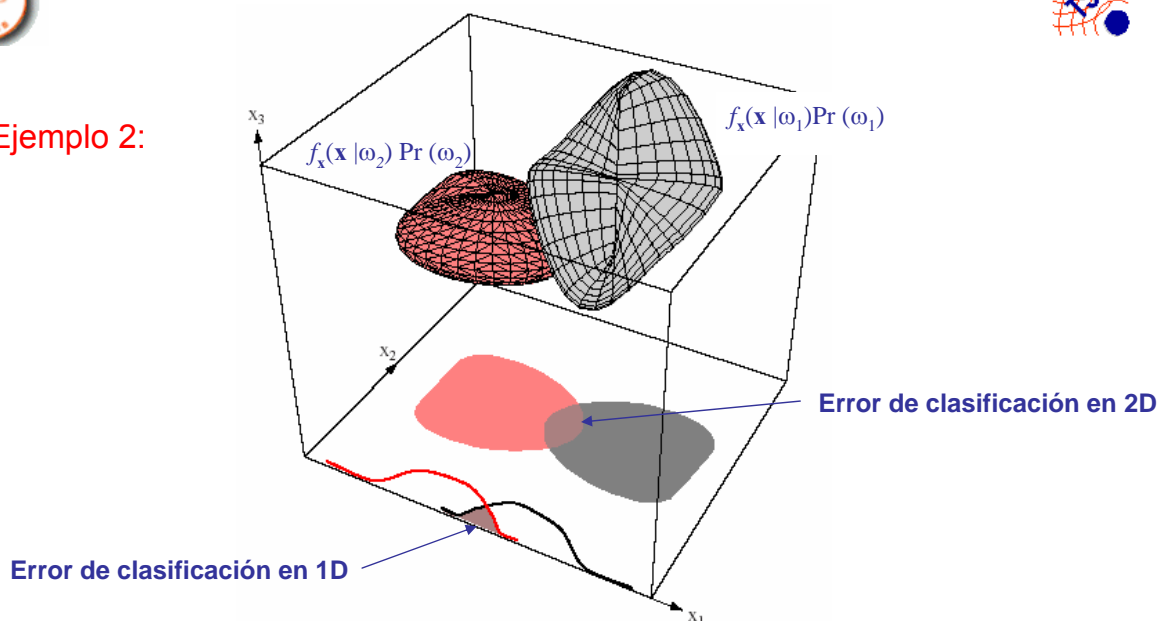
$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_d^2 \end{bmatrix} \quad r^2 = \sum_{i=1}^d \left(\frac{\mu_{i,1} - \mu_{i,2}}{\sigma_i} \right)^2$$

Quando $d \rightarrow \infty$, $r \rightarrow \infty$, $\Pr(e) \rightarrow 0$

21



Ejemplo 2:



En este caso, las pdf de dos clases están solapadas cuando se dispone de una (x_1) o dos características (x_1, x_2). Cuando se añade la tercera característica (x_3), las dos pdf aparecen completamente separadas y el error de clasificación se reduce a cero.

22



- ⇒ Sin embargo, disponer de características adicionales no siempre mejora la clasificación si se dispone de pocas observaciones para determinar parámetros adicionales.
- ⇒ Por otra parte, es posible que algunas de las características:
- No aporten información que permita clasificar mejor
 - No sean independientes de otras características
 - No se ajusten a la pdf asumida
- ⇒ Disponer de muchas características aumenta la complejidad de la solución.

La definición del número de características relevantes es necesaria en todo sistema de clasificación

23



Reducción del número de características mediante combinación lineal:

$$\mathbf{y}_k = \mathbf{W}^T \mathbf{x}_k \quad \mathbf{x} \in \mathbb{R}^d, \mathbf{y} \in \mathbb{R}^{d'}, \mathbf{W}^T \in \mathbb{R}^{d' \times d} \quad d' < d$$

Soluciones posibles para \mathbf{W} :

1. Proyectar los vectores \mathbf{x}_k de la mejor forma posible en el sentido del error cuadrático
⇒ análisis de componentes principales (PCA)
2. Proyectar los vectores \mathbf{x}_k de forma que las clases resultantes queden lo más separadas posible
⇒ análisis por múltiples discriminantes (MDA)

24



ANÁLISIS DE COMPONENTES PRINCIPALES



Disponemos de un total de N observaciones asociadas a todas las clases:

$$D = \{ \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \} \quad \mathbf{x} \in \mathbb{R}^d$$

buscamos una matriz \mathbf{W} unitaria ($\mathbf{W}^T \mathbf{W} = \mathbf{I}$) que mejor aproxime

$$\mathbf{x}_k \cong \mathbf{W} \mathbf{y}_k + \mathbf{m} \quad \mathbf{m} = \frac{1}{N} \sum_{\mathbf{x} \in D} \mathbf{x}$$

en el sentido del error cuadrático medio

$$J = \sum_{k=1}^N \|\mathbf{W} \mathbf{y}_k + \mathbf{m} - \mathbf{x}_k\|_2^2 \quad (1)$$

25



Los vectores de dimensión reducida vienen dados por:

$$\mathbf{y}_k = \mathbf{W}^T (\mathbf{x}_k - \mathbf{m}) \quad (2)$$

y la matriz de transformación es $\mathbf{W} = [\mathbf{w}_1 \quad \mathbf{w}_2 \quad \dots \quad \mathbf{w}_{d'}]$

donde los vectores \mathbf{w}_j cumplen:

$$\mathbf{S} \mathbf{w}_j = \lambda_j \mathbf{w}_j$$

$$\mathbf{S} = \sum_{\mathbf{x} \in D_i} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T$$

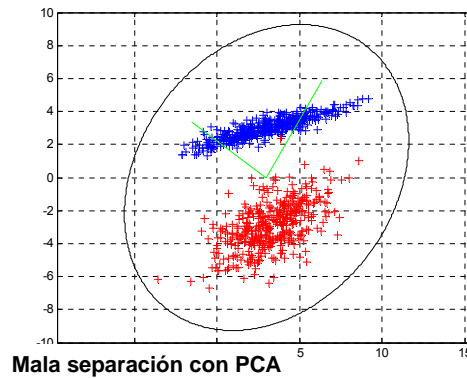
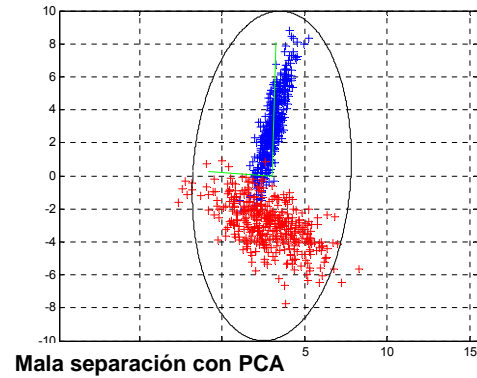
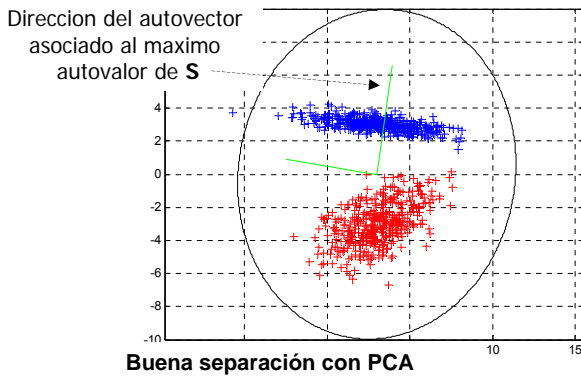
La minimización del error cuadrático medio implica utilizar en \mathbf{W} los autovectores asociados a los mayores autovalores.

**** Demostradlo, minimizando (1) con la restricción $\mathbf{w}_j^T \mathbf{w}_j = 1$ ****

26



Los autovectores definen las direcciones principales de un hiperelipsoide. Son ortogonales y apuntan en las direcciones de máxima dispersión.



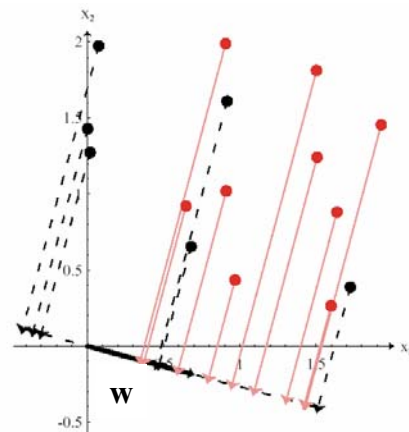
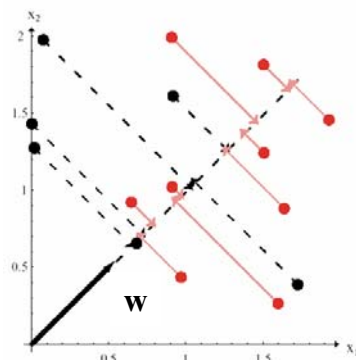
ANÁLISIS POR MÚLTIPLES DISCRIMINANTES



La transformación por componentes principales no es siempre útil para discriminar entre distintas clases. Sería mejor definir una transformación que

- Aumente la distancia inter-clases y
- Disminuya la dispersión intra-clase.

Ejemplo 1:





Construiremos una transformación \mathbf{W}^T desde un espacio de dimension d (tamaño de los vectores observados \mathbf{x}) a un espacio de dimension $d' \leq c-1$ (donde c es el numero de clases).

Necesitamos una medida de la **distancia inter-clases** y una medida de la **dispersión intra-clase**, para lo cual definimos:

$$\mathbf{m}_i = \frac{1}{N_i} \sum_{\mathbf{x} \in D_i} \mathbf{x} \quad \text{Media de los datos de la clase } i$$

$$\mathbf{m} = \frac{1}{N} \sum_{\mathbf{x} \in \{D_1, \dots, D_c\}} \mathbf{x} = \frac{1}{N} \sum_{i=1}^c N_i \mathbf{m}_i \quad \text{Media de todos los datos}$$

$$\mathbf{S}_T = \sum_{\mathbf{x} \in \{D_1, \dots, D_c\}} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T \quad \text{Dispersión total de los datos}$$

29

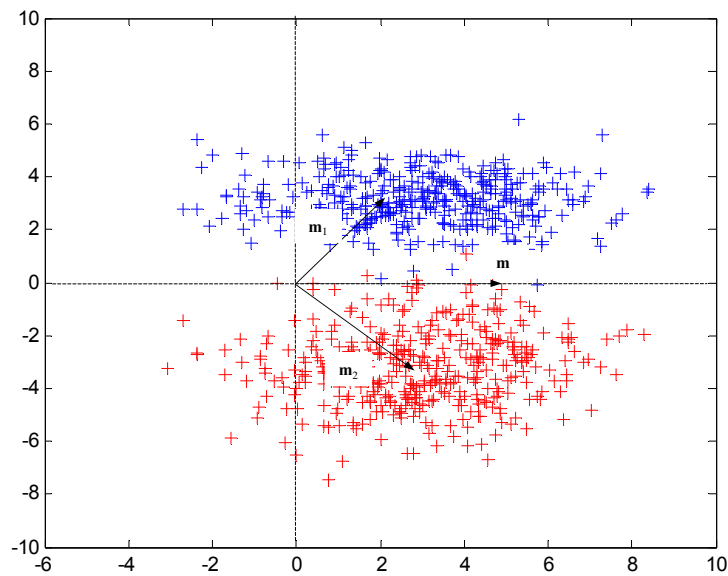


$$\begin{aligned} \mathbf{S}_T &= \sum_{\mathbf{x} \in \{D_1, \dots, D_c\}} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T = \\ &= \sum_{i=1}^c \sum_{\mathbf{x} \in D_i} (\mathbf{x} - \mathbf{m}_i + \mathbf{m}_i - \mathbf{m})(\mathbf{x} - \mathbf{m}_i + \mathbf{m}_i - \mathbf{m})^T = \\ &= \sum_{i=1}^c \sum_{\mathbf{x} \in D_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^T + \sum_{i=1}^c \sum_{\mathbf{x} \in D_i} (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T = \\ &= \sum_{i=1}^c \mathbf{S}_{C,i} + \sum_{i=1}^c N_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T = \mathbf{S}_C + \mathbf{S}_B \end{aligned}$$

Suma de matrices de dispersión intra-clases

Matriz de distancia inter-clases

30



31



La transformación a aplicar será:

$$\mathbf{y}_k = \mathbf{W}^T \mathbf{x}_k \quad \mathbf{x} \in \mathbb{R}^d, \mathbf{y} \in \mathbb{R}^{d'}, \mathbf{W}^T \in \mathbb{R}^{d' \times d} \quad d' < d$$

Las matrices de dispersión intra-clase e inter-classes quedan modificadas:

$$\mathbf{S}_C \rightarrow \mathbf{W}^T \mathbf{S}_C \mathbf{W}$$

$$\mathbf{S}_B \rightarrow \mathbf{W}^T \mathbf{S}_B \mathbf{W}$$

Si suponemos que \mathbf{y} es gaussiano, el volumen de las pdf asociadas a cada clase es proporcional a $|\mathbf{W}^T \mathbf{S}_C \mathbf{W}|^{1/2}$ y la distancia entre clases es proporcional a $|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|^{1/2}$

32



Criterio discriminante:

$$\mathbf{W} = \arg \max_{\mathbf{w}} \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_C \mathbf{W}|} \quad (1)$$

La solución para las columnas de \mathbf{W} son los $c-1$ autovectores generalizados:

$$\mathbf{S}_B \mathbf{w}_j = \lambda_j \mathbf{S}_C \mathbf{w}_j \quad (2)$$

asociados a los $c-1$ autovalores generalizados mayores.

33



La **dimensión máxima de los vectores** y es $c-1$, ya que:

1. El rango de \mathbf{S}_B es $c-1$
2. Por tanto solo existen $c-1$ autovalores distintos de cero
3. La maximización de (1) implica que solo los autovectores asociados a los autovalores distintos de cero han de incluirse en la solución.

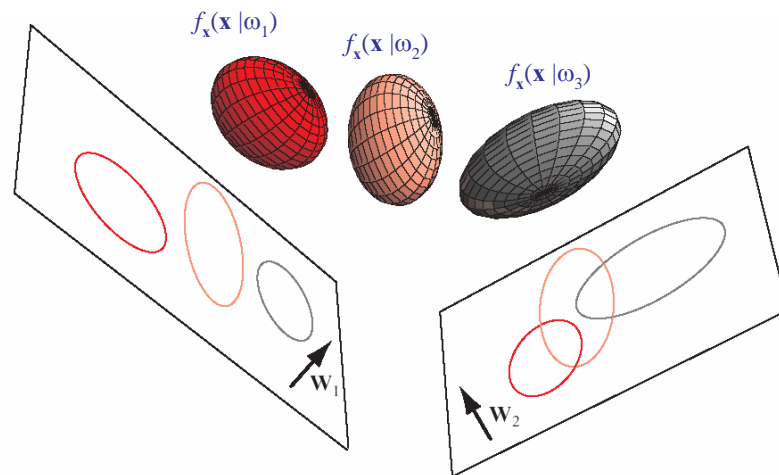
***** Demostrad que la ecuación (2) maximiza (1) *****

***** substituyendo en (1) y usando $|\mathbf{AB}| = |\mathbf{A}||\mathbf{B}|$ *****

34



Ejemplo 2:



Las distribuciones tridimensionales asociadas a tres clases distintas se proyectan sobre espacios de dimensión inferior (planos en este caso). La proyección sobre el plano definido por \mathbf{W}_1 proporciona mayor separación entre clases que el plano definido por \mathbf{W}_2 .

35



Regiones de decisión sobre los datos transformados

El cálculo de las regiones de decisión puede hacerse siguiendo las pautas de un detector bayesiano óptimo (ver tema 2.1) suponiendo gaussianidad para \mathbf{y} , con parámetros:

$$\boldsymbol{\mu}_i = \mathbf{W}^T (\mathbf{m}_i - \mathbf{m}) \quad i = 1, \dots, c$$

$$\mathbf{C}_i = \mathbf{W}^T \mathbf{S}_{c,i} \mathbf{W}$$

36



2.2.5 CONCLUSIONES



- Si se puede suponer una forma paramétrica para $f_{\mathbf{x}}(\mathbf{x}|\omega_i)$ entonces la fase de entrenamiento del clasificador se reduce a la estimación de los parámetros
- Pueden utilizarse dos soluciones para la estimación de parámetros: ML (más simple computacionalmente) o bayesiana (si se dispone de conocimiento a priori sobre los parámetros)
- Los problemas derivados del exceso de dimensionalidad pueden reducirse mediante el uso de múltiples discriminantes